



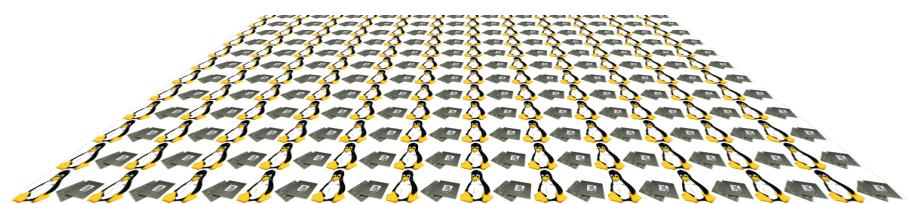
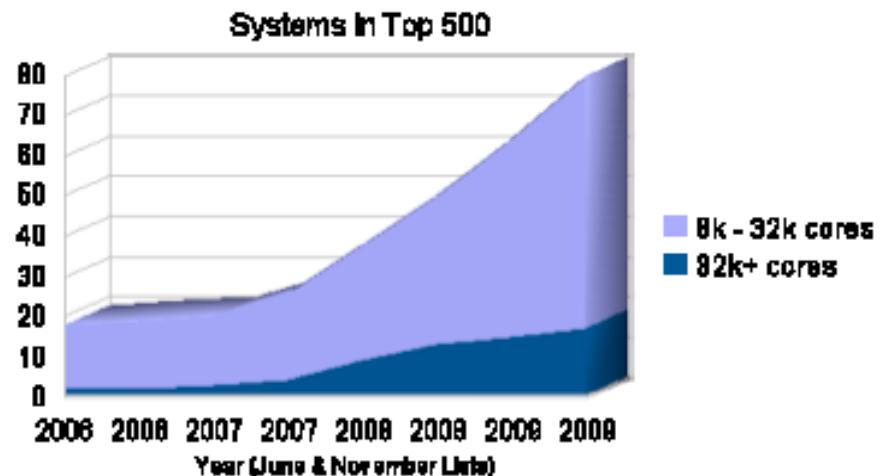
# Debugging Applications in Jaguar XT5

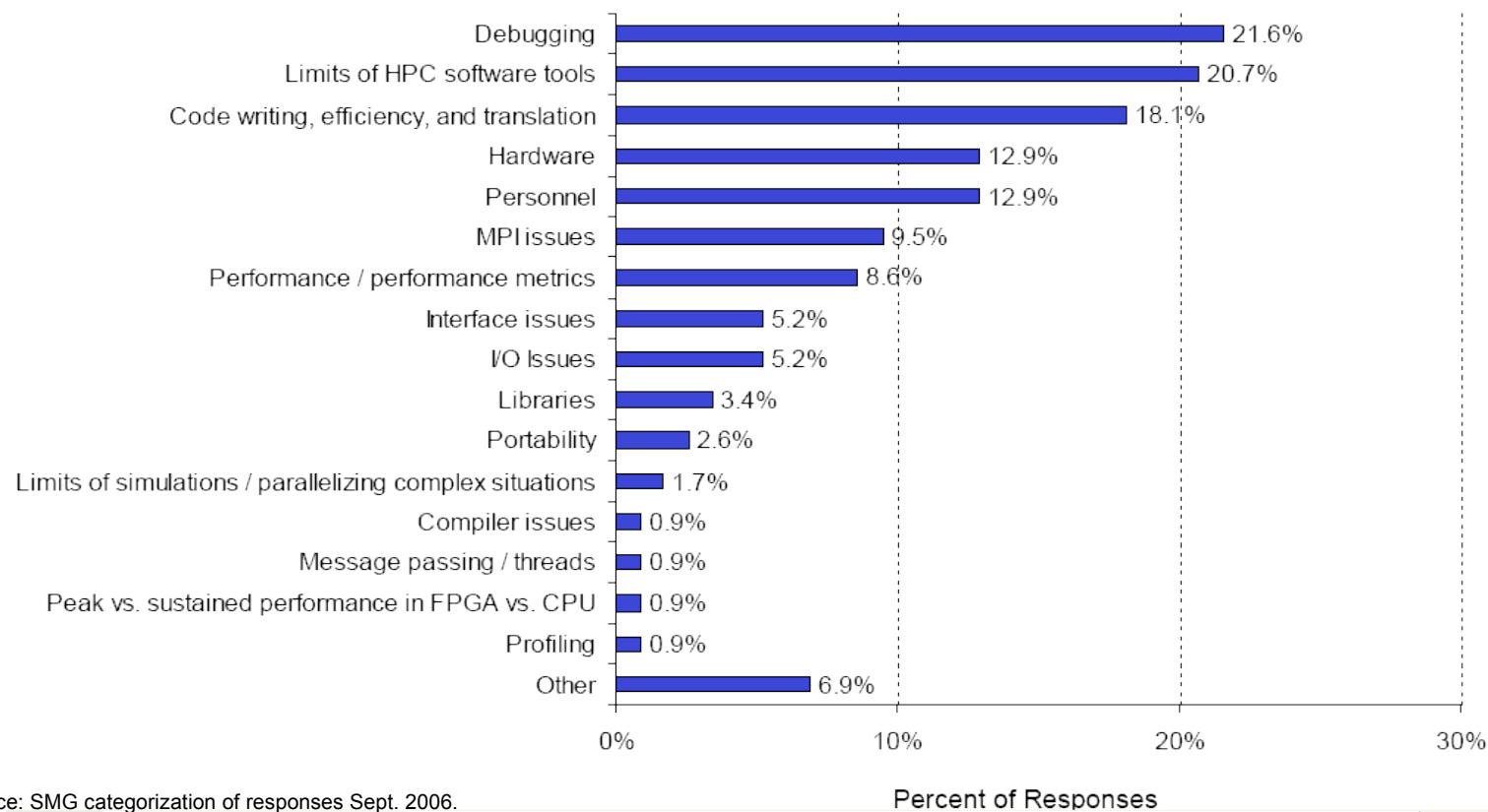
Presented by:

Oscar Hernandez, Rich Graham  
ORNL

David Lecomber, David Maples  
Allinea

- Processor counts growing rapidly
- GPUs entering HPC
- Large hybrid systems imminent
- But what happens when software doesn't work?





Source: SMG categorization of responses Sept. 2006.

Percent of Responses

- Shortening application development time is vital
- Debugging is now the major issue !

- HPC tools company since 2001
  - DDT - Debugger for MPI, threaded/OpenMP and scalar
  - OPT - Optimizing and profiling tool for MPI and non-MPI
  - DDTLite – Parallel Debugging Plugin for Microsoft Visual Studio 2008 SP1 and above
- Large European and US customer base
  - Ease of use – means tools get used
  - Users debugging regularly at all scales
  - Scalable interface – easy to use at 1 or 100,000s of cores
- Looking to the future
  - In use at Petascale, collaboration with ORNL
  - Building debuggers that scale on Jaguar XT5
  - GPU product in Beta

- We need good debuggers
  - Instrument the code: printf - is not a time-effective solution
  - Command line debugging: doesn't scale
  - Complexity scales..
    - User overload: beware the deluge of information
    - Good tools can present information better
    - Aggregation of control across cores
  - Good software tools lead to
    - Faster development
    - More effective maintenance/deployment

- A powerful and highly intuitive tool
  - Traditional focus has been HPC
- Cross-platform support
  - Linux, Solaris, AIX, Super UX, Blue Gene O/S
  - Blue Gene, Cell, x86-64, ia64, PowerPC, Sparc, NEC, NVIDIA
  - GNU, Absoft, IBM, Intel, PGI, Pathscale, Sun compilers
- Across all MPI and OpenMP implementations
  - From low end to high end
- Support for all scheduling systems
  - SGE, PBS, LSF, MOAB, ...
  - Flexible, powerful, easy to use queue submission

- A sophisticated GUI aids development
  - Helps the user to control parallel execution
  - Helps to find and focus on potential problems
  - User controls actions.....
    - Set breakpoints, watchpoints, lock step, align stacks etc
    - Can focus in on individual threads/processes
    - Can quickly and easily create groups of processes
      - Controlled in the same way as individual
      - See stack across all processes
    - All based on code logic !

The screenshot shows the Allinea Distributed Debugging Tool (DDT) interface. At the top, there's a menu bar with Session, Select, Search, View, Code, and Help. Below the menu is a toolbar with various icons. The main window has a title bar "Allinea Distributed Debugging Tool - /home/matt/ddt/examples/hello.c". On the left, there's a "Current Group" dropdown set to "Crash" and a "Threads" list showing 15 threads (0-14). Below that are sections for "Workers" (1-15), "Crash" (3, 5, 8), and "Root" (0). On the right, there's a "Project Files" tree with "Header Files" and "Source Files" (hello.c selected). The code editor shows the following C code:

```
96    printf("I have %d arguments.\n", argc);
97    printf("\tHow many did I say?\n");
98
99    printf("They are:\n");
100   for(i=0;i<argc;i++)
101     On this line:
102     printf("%s\n", argv[i]);
103   Workers (3,0, 5,0, 8,0)
104   All (3,0, 5,0, 8,0)
105   Crash (3,0, 5,0, 8,0)
106   e:\n");
107   Breakpoint for Crash
108   printf("%s\n", *environ);
```

Loading DDT in Jaguar XT5

```
%module load ddt
```

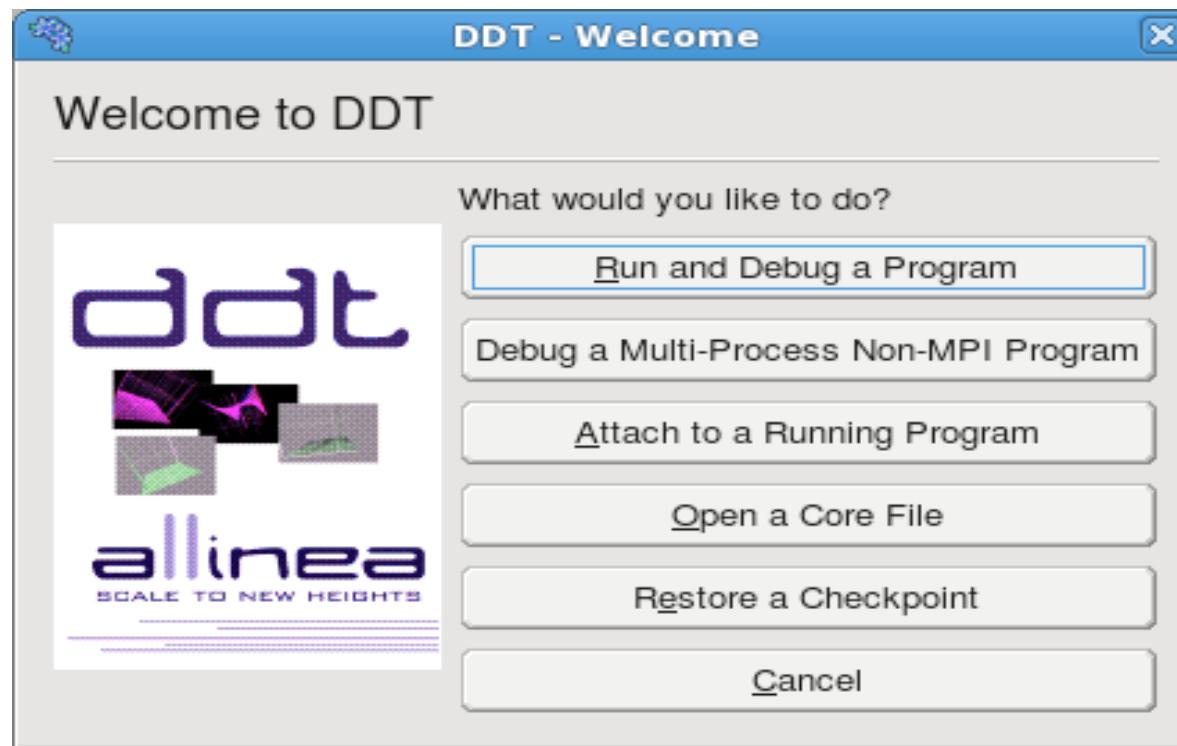
Compiling your application with debugging option

```
%cc -g myapplication.c -o myapplication
```

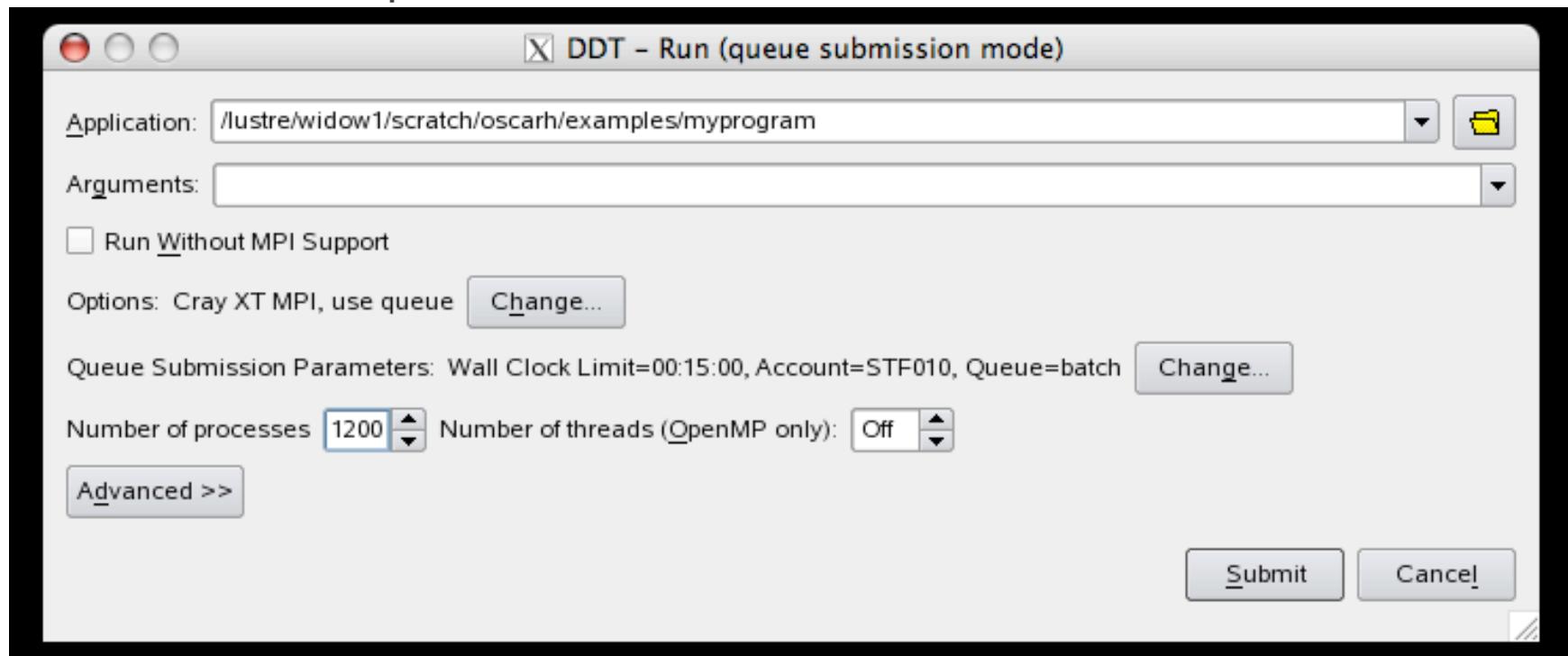
Invoking DDT

```
%ddt myapplication
```

- Ability to run in numerous modes

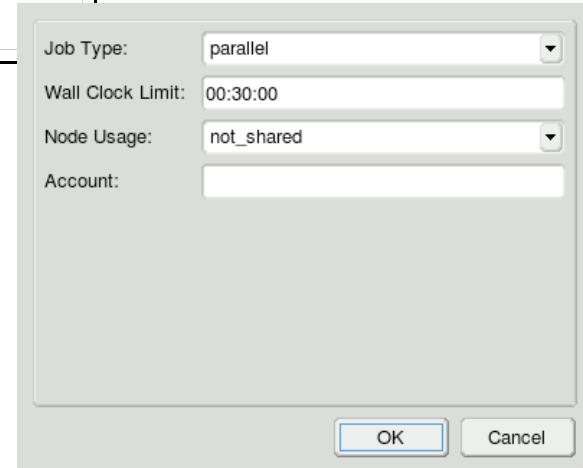


- DDT has a GUI to submit jobs
- Can customize submission scripts with GUI
- You can specify MPI versions, input data information, advanced submission options.



- Enhanced Queuing Submission
  - Simple mechanism to customise options
  - Submit non-mpi jobs to a queue
  - Set pre-defined variables in submission script
    - System administrator configures

```
# Name: LoadLeveler
# JOB_TYPE_TAG: {type=select,options=parallel|serial,label="Job Type",default=parallel}
# WALL_CLOCK_LIMIT_TAG: {type=text,label="Wall Clock Limit",default="00:30:00",mask="09:09:09"}
# NODE_USAGE_TAG: {type=select,options=not_shared|shared,label="Node Usage",default=not_shared}
# ACCOUNT_TAG: {type=text,label="Account",global}
```



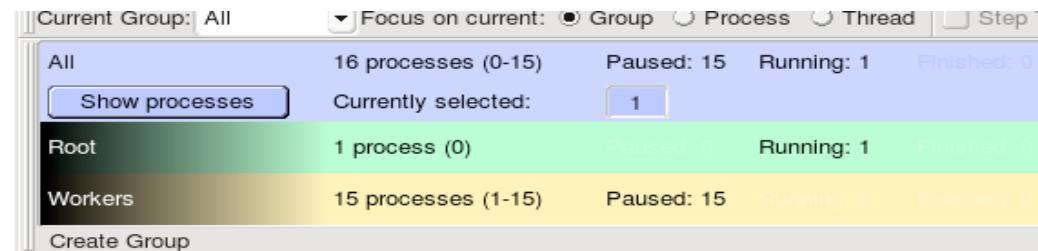
- Alter the pre-defined variables at run-time
  - User selects own specific values

- Scalar Features

- Advanced C++ support including STL, namespaces, virtual functions and templates
- Advanced Fortran 90, 95 and 2003 support including modules, allocatable data, pointers and derived types

- Multithreading & OpenMP features

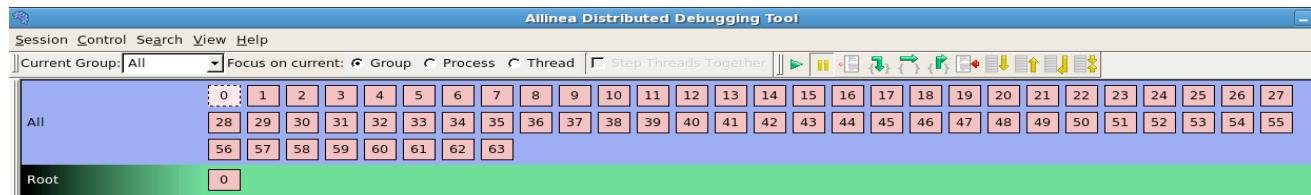
- Perform actions individually or collectively



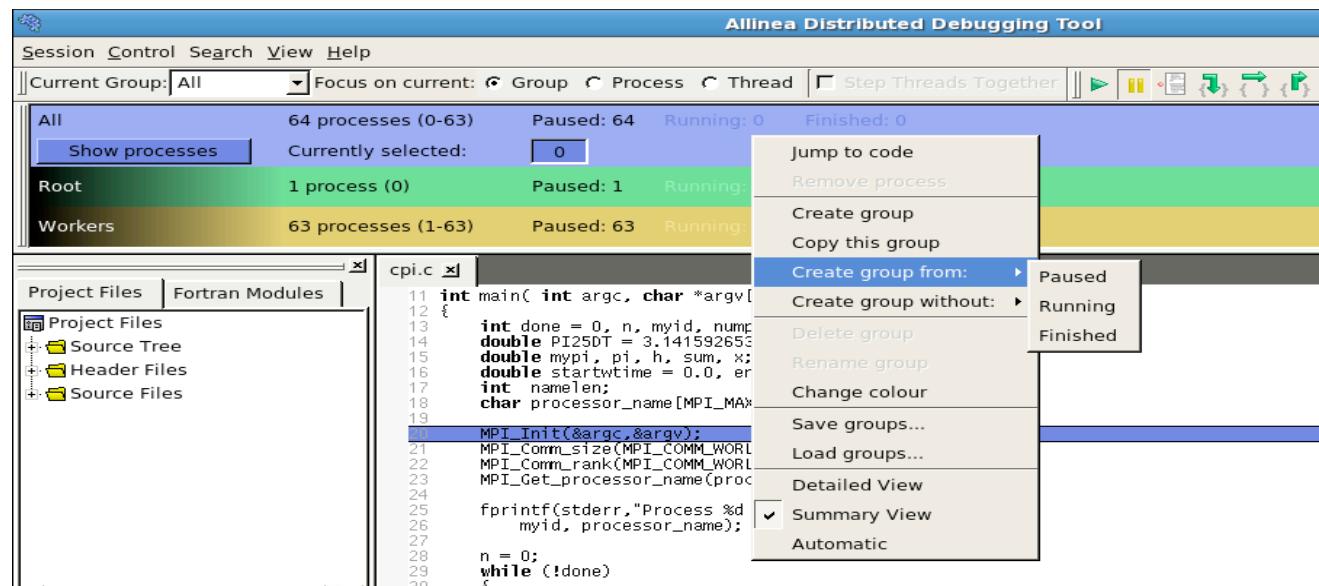
- MPI features

- Control processes
- Individually or by groups

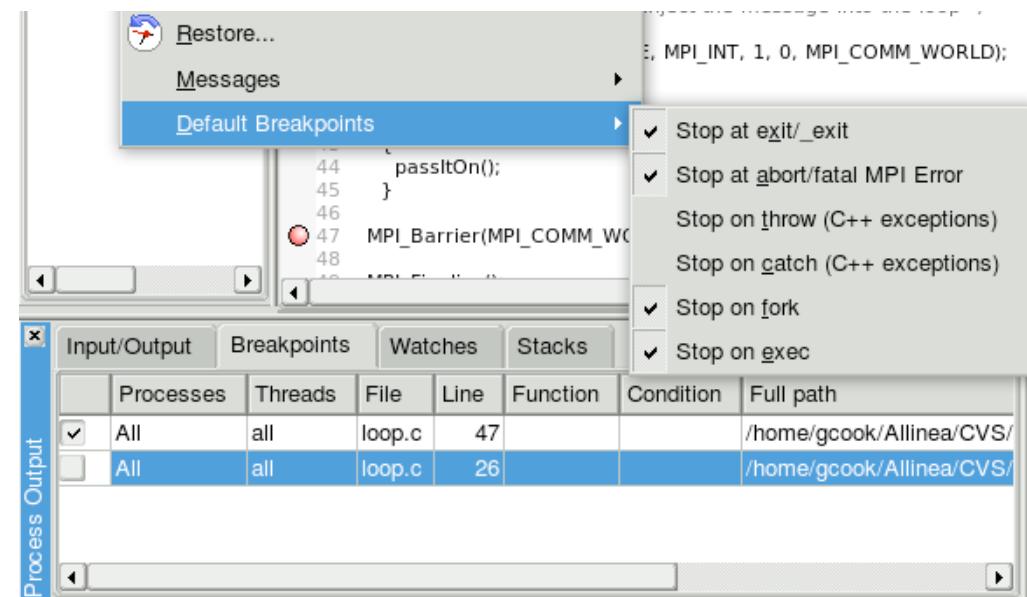
- Process View (Summary or Detail)
  - Chosen automatically or user driven
  - Create new groups based on process “status”



- Process View (Summary or Detail)
  - Chosen automatically or user driven
  - Create new groups based on process “status”
  - Control 1000s of processes in the same way
  - All based on code logic



- Quickly and easily control your debug session
  - Breakpoints
    - Default settings (enable/disable)
    - Double click to add them
    - For all processes, groups or a single process
    - Enable/disable from window panel
    - Add a condition



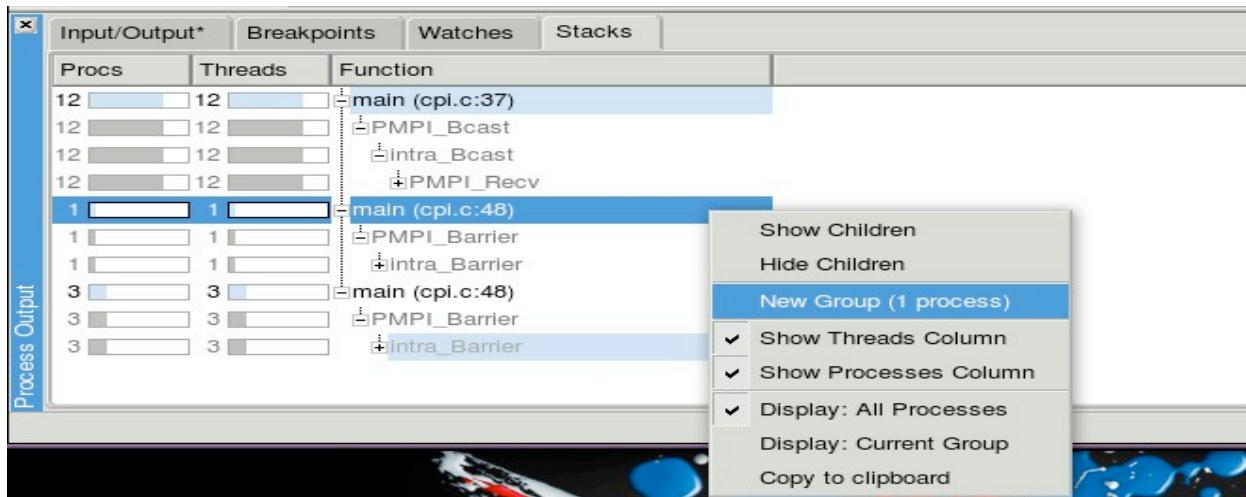
- Quickly and easily control your debug session

The screenshot shows the DDT (Debugging and Profiling Tool) interface. At the top is a code editor displaying C code for MPI calculations. Below it is a 'Watches' window with tabs for Input/Output\*, Breakpoints, Watches, and Stacks. The Watches tab is selected, showing a table of variables with their current values. A context menu is open over the variable table, with 'Add To Watches' highlighted. The menu also includes options like View As, View Array (MDA), View Across Processes (CPC), View Across Threads (CTC), View As Vector (C/C++ only), Get Address, Dereference Pointer, Check Pointer Is Valid, View Pointer Details, Find Variable In Files, Evaluate With Wizard, Delete, and Delete All.

Expression	Value
h	0.01
myid	2
n	100
pi	6.95327155555555e-010

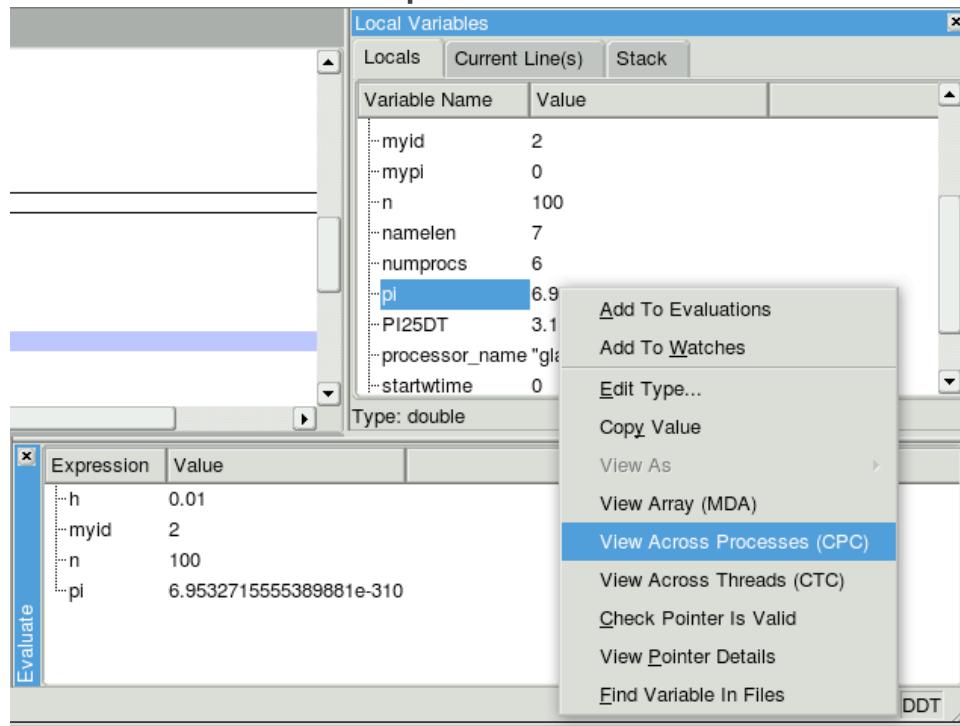
## – Watchpoints

- Select from evaluate, locals or current line
- Set if variable changes
- Program halts if change occurs
- Single process only

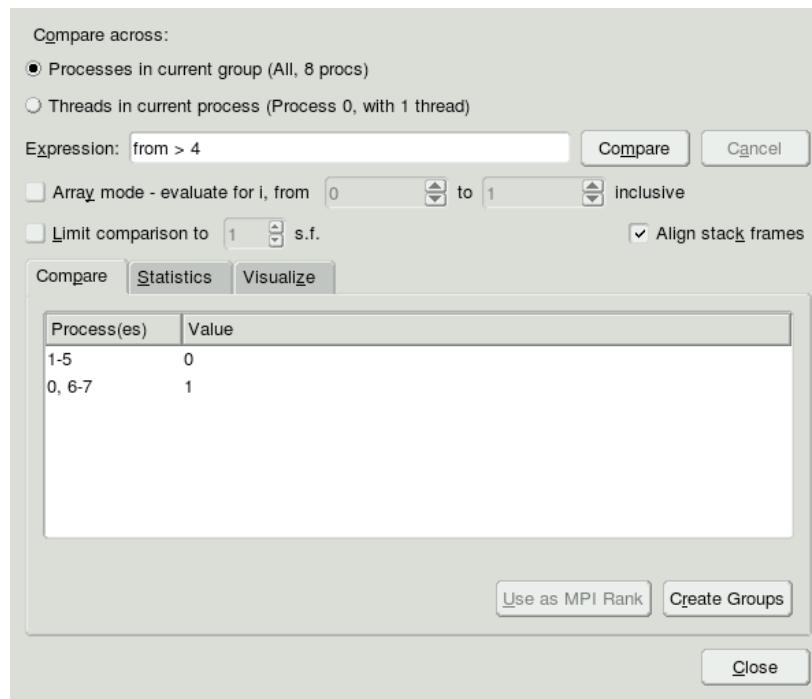


- Parallel Stack View
  - Finds rogue processes faster
  - Identifies classes of process behaviour easily
  - Allows rapid grouping of processes
  - All based on code logic

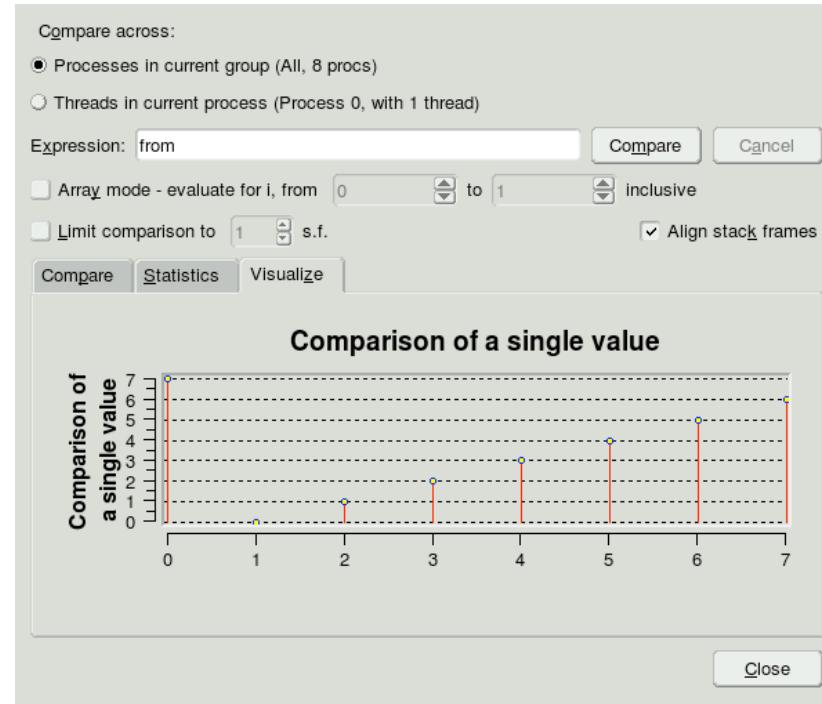
- Cross Process/Thread Comparison (CPC)
  - Run from evaluate, locals or current line
  - Variable or expression



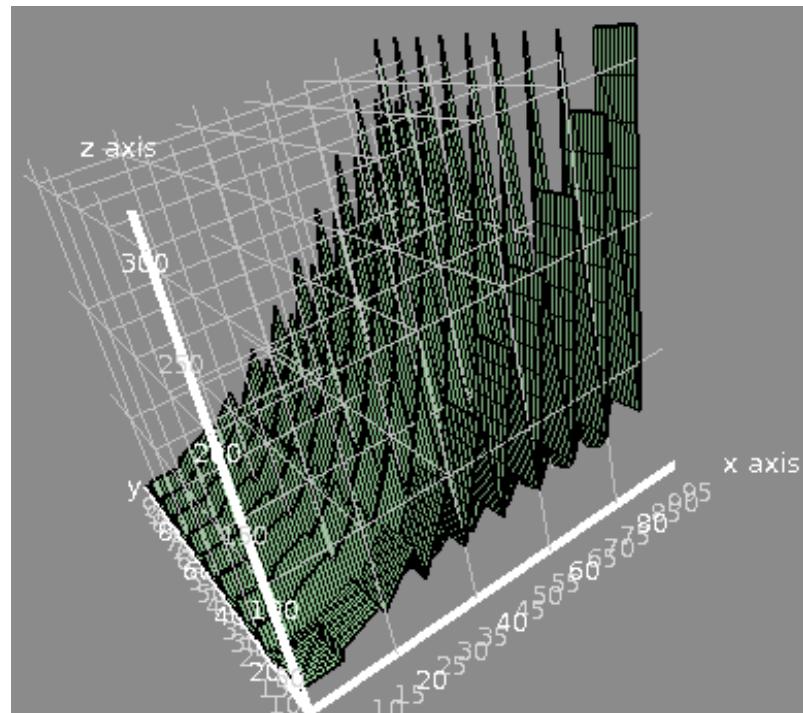
- Cross Process/Thread Comparison (CPC)
  - Run from evaluate, locals or current line
  - Variable or expression
  - Create groups from results



- Cross Process/Thread Comparison (CPC)
  - Run from evaluate, locals or current line
  - Variable or expression
  - Create groups from results
  - Visualize

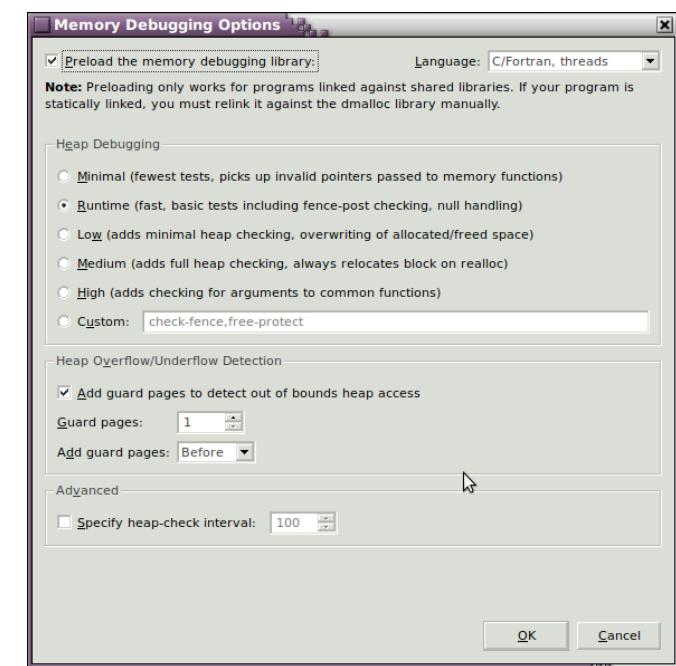
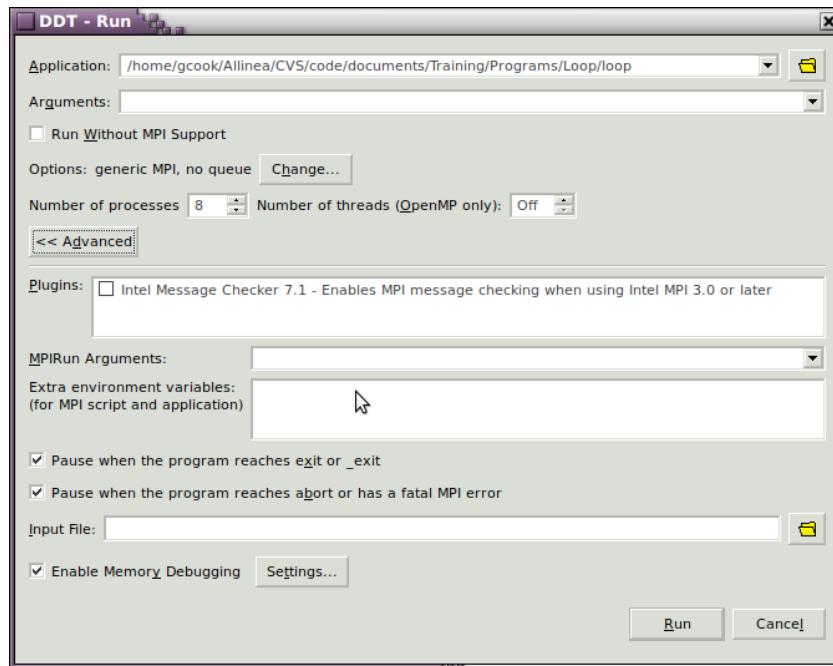


- Cross Process/Thread Comparison (CPC)
  - Run from evaluate, locals or current line
  - Variable or expression
  - Create groups from results
  - Visualize
  
- Multi Dimensional Array viewer
  - 3D OpenGL array viewer
  - Export csv file for post analysis



- **Simple to use**

- Tick check box and use defaults.....
- .....or select “Settings”
- Enhanced heap debugging settings
- Enable/disable guard pages (before or after)



Check your overall memory usage

The screenshot shows the Allinea Distributed Debugging Tool (DDT) interface. On the left, there's a file browser with project files like Header Files, Source Files, and various Fortran source files (check.f90, dlanv.f, gendata.f90, matnrm.f90, mod\_trsol.f90, numroc.f, solve.f90, trsol.f90). Below it is a terminal window showing command-line output. In the center, a callout box contains the text "Check your overall memory usage". To the right, there are two bar charts under the "Overall Memory Stats" tab. The first chart, "Total bytes allocated/thread", shows a series of red bars for threads P.0 through P.7, with values decreasing from approximately 850,000 to 450,000. The second chart, "total number of memory allocation/deallocation calls", shows red bars for allocation calls and green bars for deallocation calls, with both generally decreasing from P.0 to P.7.

Overall Memory Stats

Total bytes allocated/thread

Thread	Total Allocated Bytes
P.0	~850,000
P.1	~480,000
P.2	~450,000
P.3	~450,000
P.4	~450,000
P.5	~450,000
P.6	~450,000
P.7	~450,000

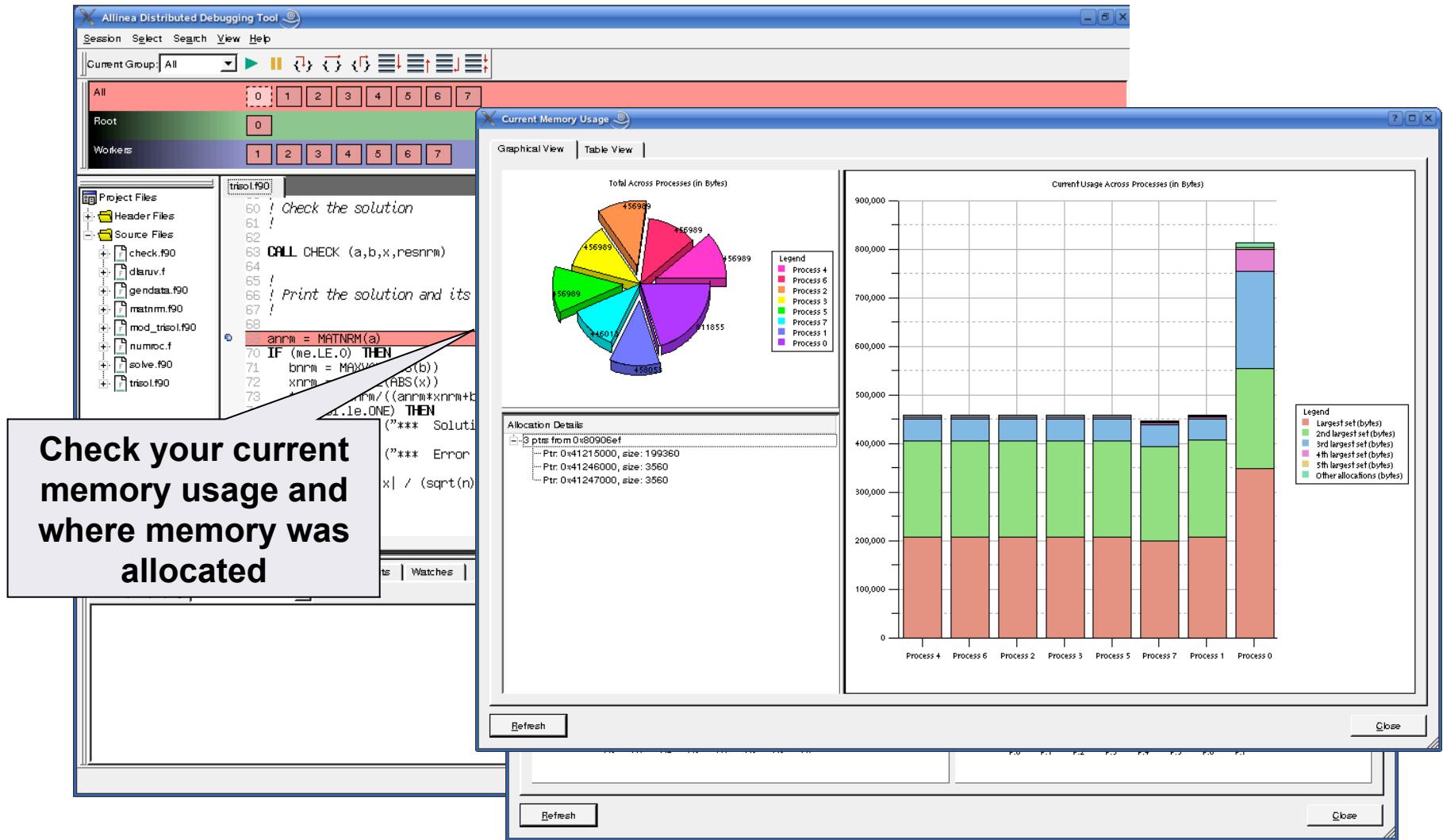
Legend: total allocated bytes (red), total freed bytes (green)

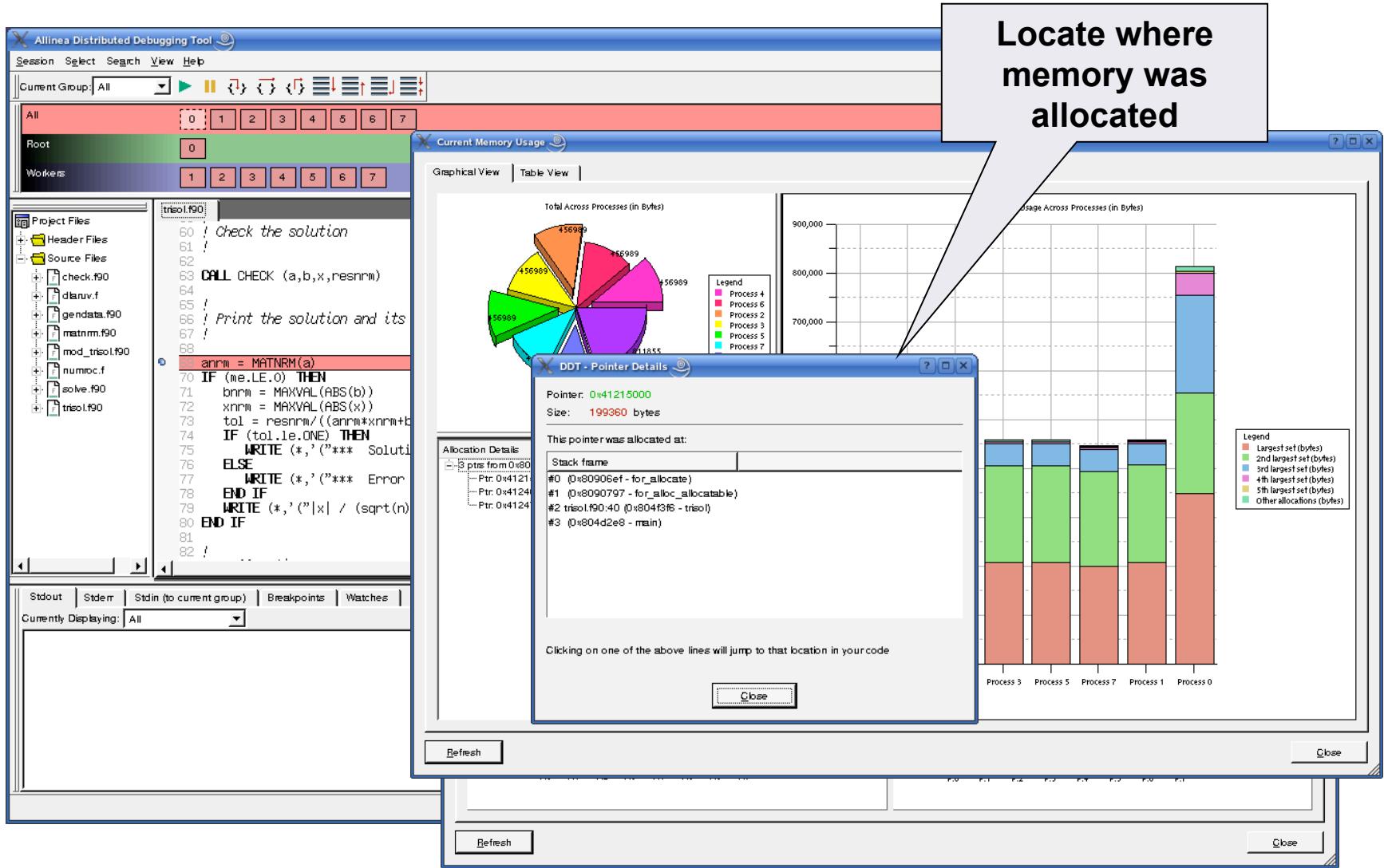
total number of memory allocation/deallocation calls

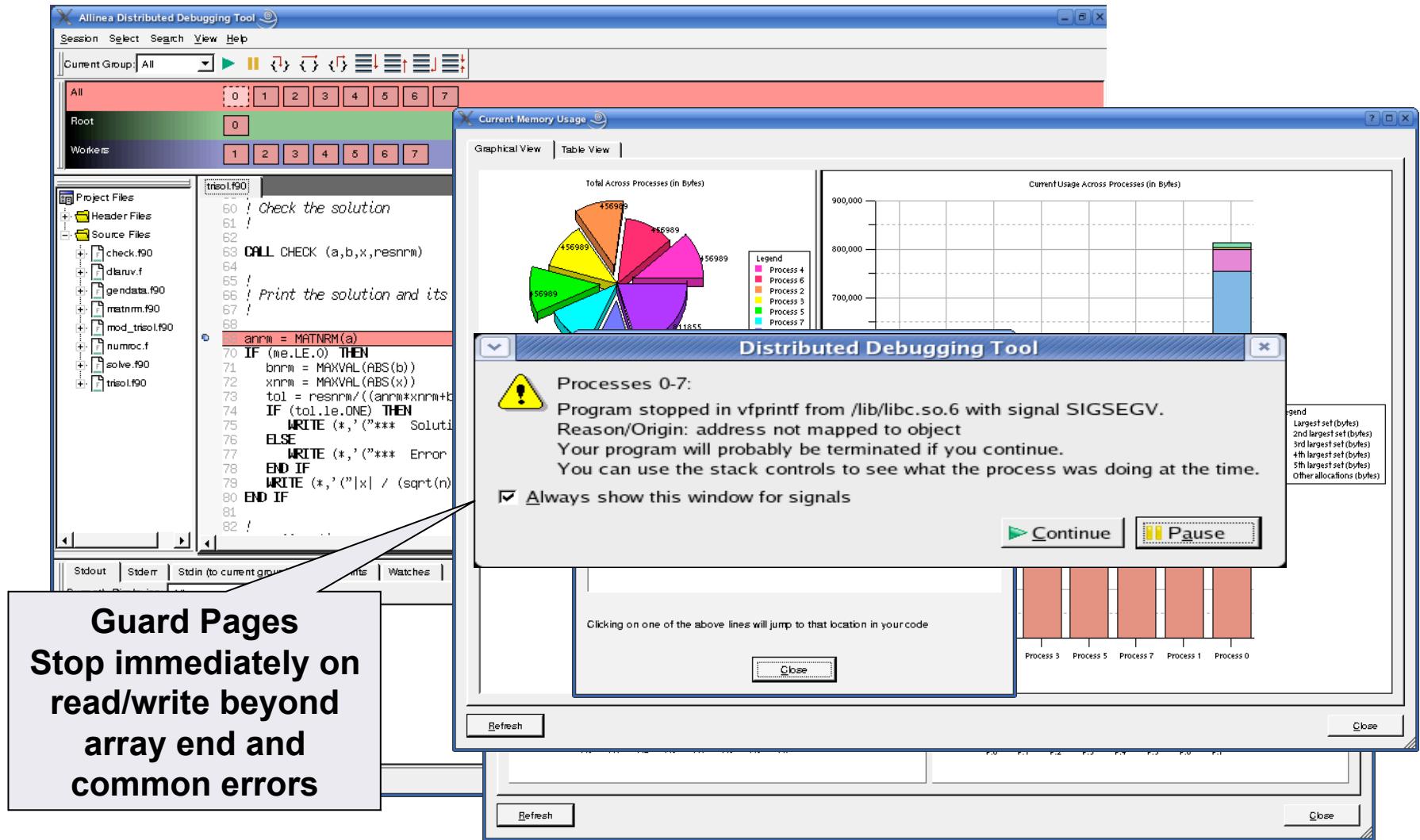
Thread	Allocation Calls	Deallocation Calls
P.0	~200	~100
P.1	~120	~25
P.2	~120	~25
P.3	~120	~25
P.4	~120	~25
P.5	~120	~25
P.6	~120	~25
P.7	~120	~25

Legend: allocation calls (red), deallocation calls (green)

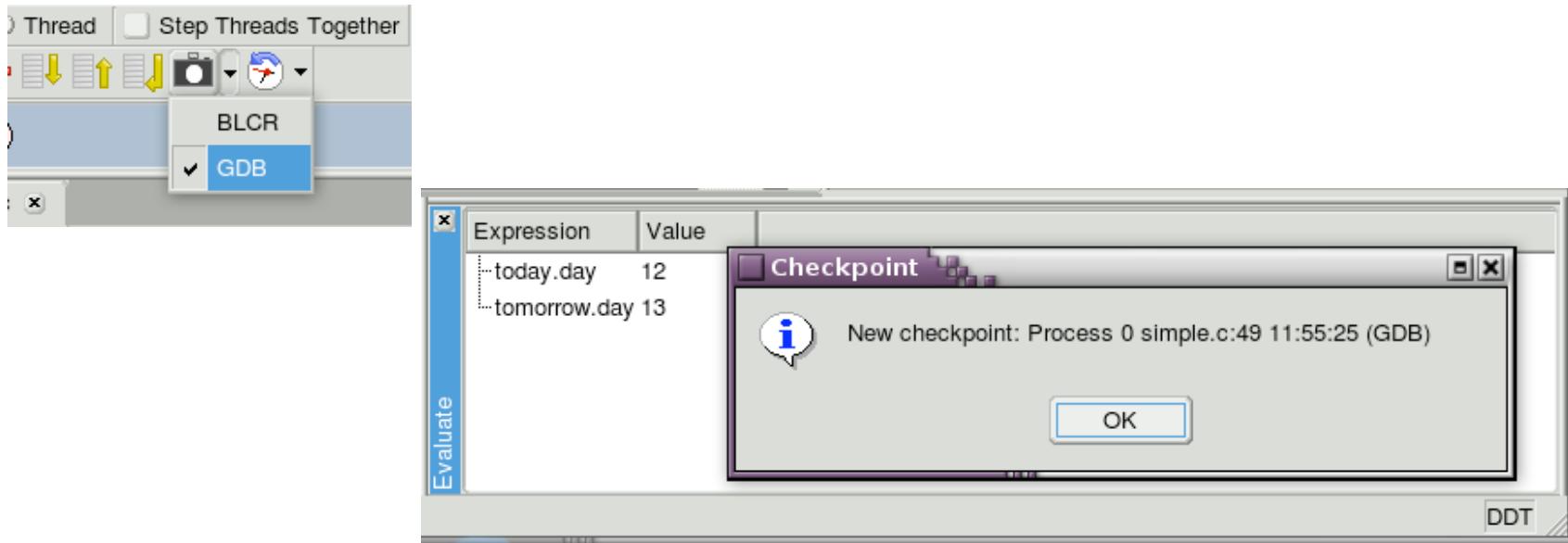
Refresh Close





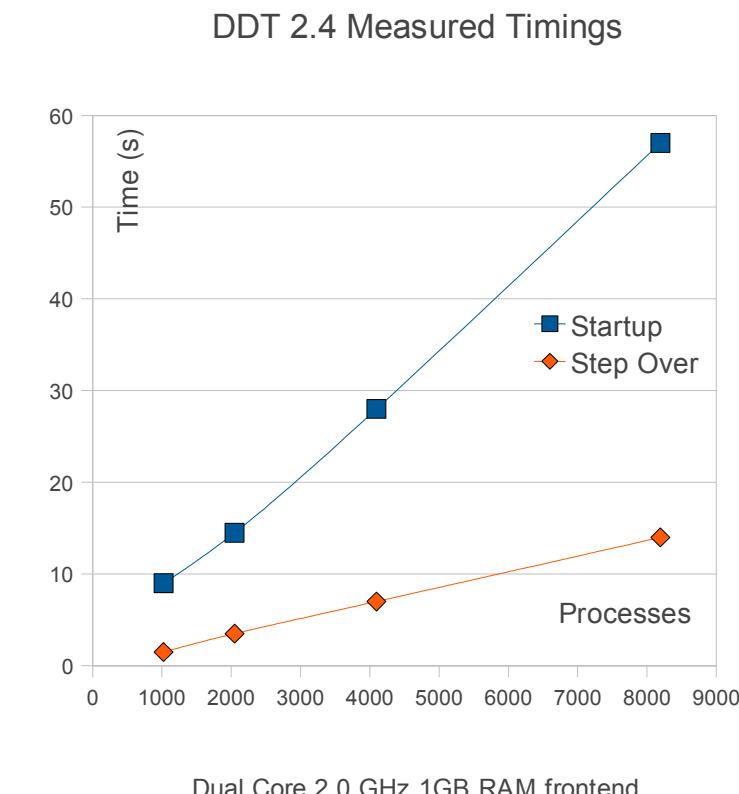


- GDB
  - Single session, no option to save/restore
  - MPI support
  - No thread support
  - Still in the experimental stage

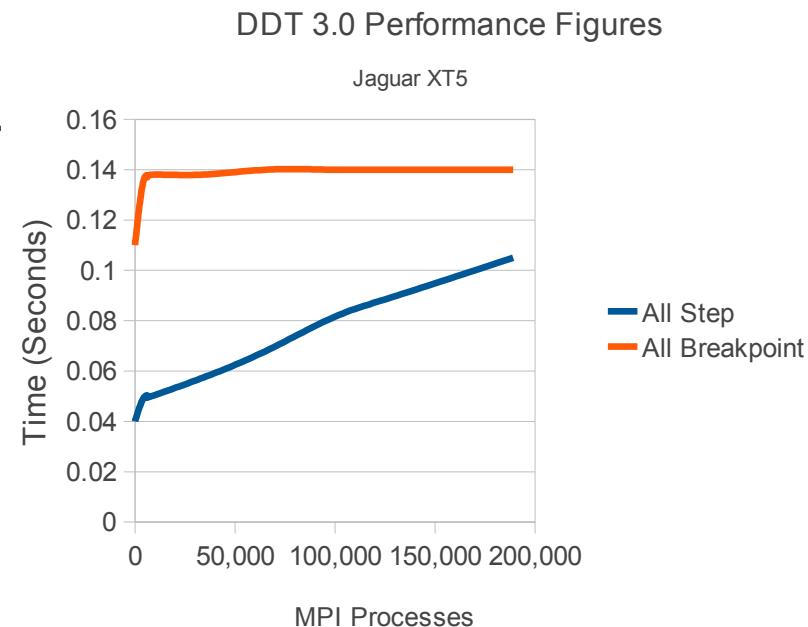


- What has been achieved in last 12 months?
- Scalable GUI
  - For the first time debug 5,000 with same ease as 100
  - At a glance full stack and status of all processes
- 10x improvement in scale limits
  - Iterative improvement has brought benefit
  - Debugging 8k processes is comfortable
  - Regular users at 4k cores
  - Test rig emulating 32k cores at native speed
- High end platforms
  - BlueGene/P support added to list Q2/09
  - Cray XT4, XT5 users at scale - 8k
  - Ranger at TACC – Infiniband Sun Constellation cluster - 8k

- Changing DDT communication network to be a shallow tree
  - Current limits ~4-16k processes
    - Linear performance
  - Progress towards >> 32k processes
    - Target of 232k



- Changing DDT communication network to be a shallow tree
  - Current limits ~4-16k processes
    - Linear performance
  - Progress towards >> 32k processes
    - Target of 232k
  - Massive increase in scalability
    - Aggregation of control messages..
    - ....and responses
  - Logarithmic performance



# Presenting Data, Usefully

The screenshot shows two windows from the allinea debugger. The top window is a variable viewer titled 'Current Line(s)' showing the values of variables 'my\_rank' (0) and 'p' (16). A tooltip '16/16 processes equal' is visible over the 'p' entry. The bottom window is a histogram comparison dialog with the expression 'my\_rank % 3'. It displays a table of process counts for each value (0, 1, 2) across 189120 processes. The table data is as follows:

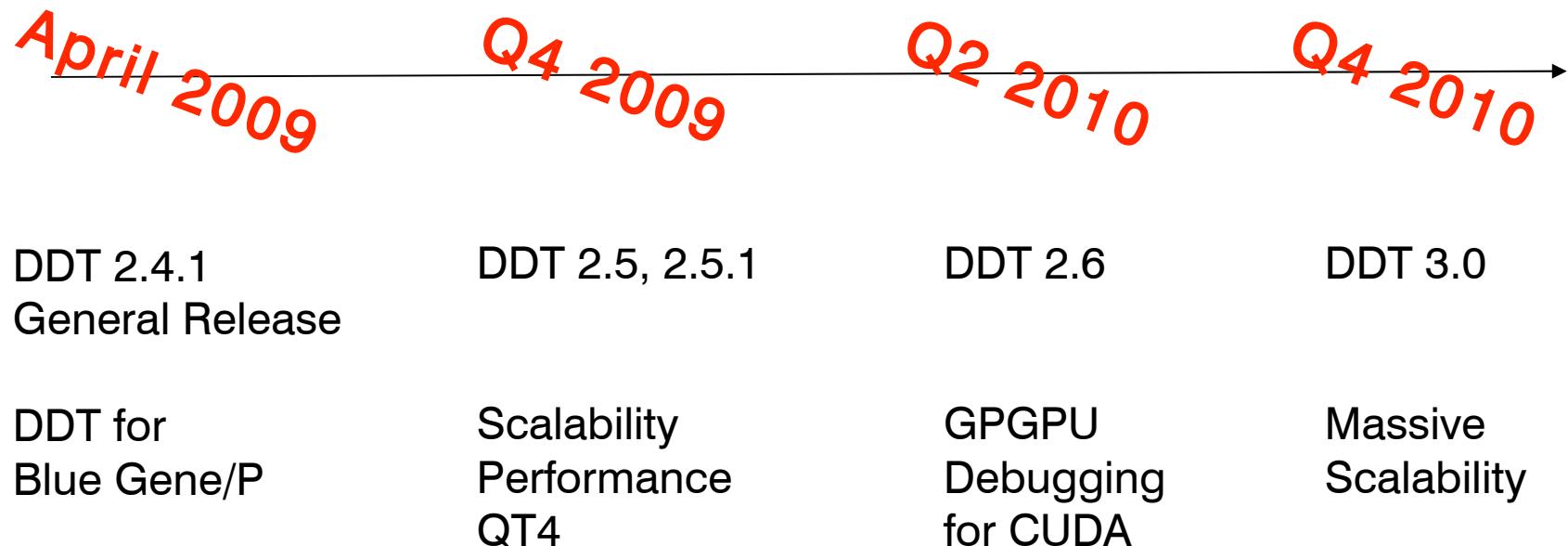
Value	Process(es)
0	0,3,6,9,12,15,18,21,24,27,30,...
1	1,4,7,10,13,16,19,22,25,28,31,...
2	2,5,8,11,14,17,20,23,26,29,32,...

On the right side of the histogram dialog, there are numerical summary statistics:

- Count: 189120
- Filtered: 0
- Errors: 0
- Aggregate: 0
- Numerical: 189120
- Sum: 189120
- Minimum: 0
- Maximum: 2

- Gather from every node
  - Potentially costly – if all data different
  - Easy if data mostly same
  - New ideas
    - Aggregated statistics
    - Probabilistic algorithms optimize performance – even in pathological case
  - ~130ms for 130,000 cores
- Watch this space!
  - With a fast and scalable architecture, new things become possible

- DDT is the first Petascale debugger..
  - A debugging tool has finally caught up with the hardware!
    - Work is in progress to port every feature for scale
    - Memory debugging, data visualization, ....
  - How can the infrastructure be built upon?
    - Does DDT offer the right framework for collaboration?
    - Can we encourage a codebase of user-generated MPI tools/utilities?
- ... but large clusters are a fraction of HPC
  - Most parallel development starts smaller
  - Is now starting even smaller: GPUs



# Questions & Answers